# *Classifying Matrices with a Spectral Regularization*

Ryota Tomioka & Kazuyuki Aihara

University of Tokyo / Fraunhofer FIRST

2007/06/23

THE UNIVERSITY OF TOKYO

FIRST

Fraunhofer Institut Rechnerarchitektur und Softwaretechnik
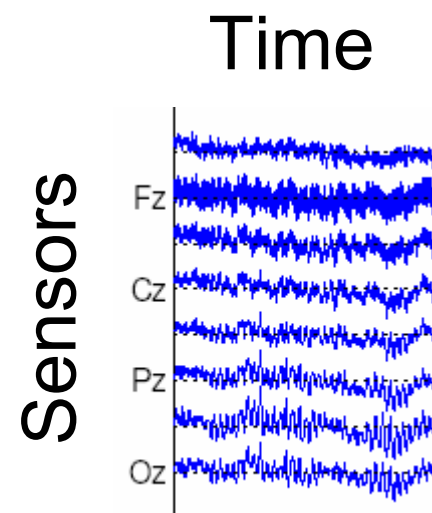
BBCI berlin brain computer interface

# Outline

- Method
  - Discriminative model that factorizes using the spectral $\ell_1$-regularization.
  - Penalized empirical loss minimization (convex!).
- Implementation
  - Dual formulation.
  - Linear Matrix Inequality.
  - Interior point method.
- Application
  - Motor-imagery EEG classification.
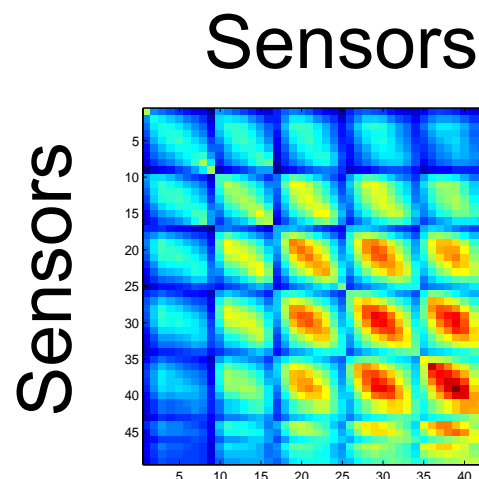- Summary

# Examples of Matrix Inputs

- Multivariate Time Series

$$X =$$

Time



- Second order statistics

$$X =$$

Sensors

# Problem Setting

The Input                          Class Label

$$X \Longrightarrow y \in \{+1, -1\}$$

$$R \times C$$

$$f(X; W, b) = \mathrm{Tr}\left[W^\top X\right] + b$$

$$(W \in \mathbb{R}^{R \times C},\ b \in \mathbb{R})$$

Spectral $\ell_1$-regularization (sum of singular-values):

$$\Omega(W) = \sum_{c=1}^{r} \sigma_c[W]$$

# Interpreting the Model

Using the singular-value decomposition:

$$W = U \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{pmatrix} V^\top = \sum_{c=1}^{r} \sigma_c \boldsymbol{u}_c \boldsymbol{v}_c^\top$$

The classifier can be written as:

$$f(X) = \mathsf{Tr}\left[\left(\sum_c \sigma_c \boldsymbol{u}_c \boldsymbol{v}_c^\top\right)^\top X\right] + b$$

$$= \sum_{c=1}^{r} \sigma_c \boldsymbol{u}_c^\top X \boldsymbol{v}_c + b$$

# Interpreting the Model

Using the singular-value decomposition:

$$W = U \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{pmatrix} V^\top = \sum_{c=1}^{r} \sigma_c \boldsymbol{u}_c \boldsymbol{v}_c^\top$$

The classifier can be written as:

$$f(X) = \text{Tr}\left[ \left( \sum_c \sigma_c \boldsymbol{u}_c \boldsymbol{v}_c^\top \right)^\top X \right] + b$$

$$= \sum_{c=1}^{r} \textcolor{red}{\sigma_c} \textcolor{blue}{\boldsymbol{u}_c^\top X \boldsymbol{v}_c} + b$$

<span style="color:red">Linear combination</span>    <span style="color:blue">Features (projected inputs)</span>

# Comments on Related Methods

LASSO:

$$\Omega_{LASSO}(W) = \sum_{(i,j)} \left| W_{ij} \right|$$

Ridge penalty:

$$\Omega_2(W) = \frac{1}{2} \sum_{i,j} W_{ij}^2 = \sum_{c=1}^{r} \sigma_c^2 [W]$$

Spectral $\ell_1$-regularization:

$$\Omega_1(W) = \sum_{c=1}^{r} \sigma_c [W]$$

# The Problem

(P) $\displaystyle \min_{\substack{W \in \mathbb{R}^{R \times C}, \\ b \in \mathbb{R}, \\ \boldsymbol{z} \in \mathbb{R}^n}} \quad \frac{1}{n} \sum_{i=1}^{n} \ell_{LR}(z_i) + \frac{\lambda}{n} \|W\|_1 ,$

<span style="color:red">Lagrange multipliers</span>

s.t. $\quad y_i \left( \mathsf{Tr}\left[ W^\top X_i \right] + b \right) = z_i \quad$ <span style="color:red">$(\alpha_i)$</span>

$$(i = 1, \ldots, n),$$

$$\ell_{LR}(z) := \log \left( 1 + \exp(-z) \right),$$

$$\|W\|_1 := \sum_{c=1}^{r} \sigma_c \left[ W \right]$$

# Implementation

- Dual Formulation

- Linear Matrix Inequality

- Interior Point Method

# The First Trick:
# The Dual Optimization Problem

(D)

$$\min_{0 \le \boldsymbol{\alpha} \le 1} \sum_{i=1}^{n} \ell^*_{\mathsf{LR}}(\alpha_i)$$

The fit must be simple (large entropy)

Residual of the fit must be small

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i y_i = 0,$$

$$\left\| \sum_{i=1}^{n} \alpha_i y_i X_i \right\|_{\infty} \le \lambda,$$

$\ell_{\infty}$-norm

$$\ell^*_{\mathsf{LR}}(\alpha) := \alpha \log \alpha + (1 - \alpha) \log (1 - \alpha),$$

$$\|X\|_{\infty} := \max_{c} \sigma_c [X].$$

# The Second Trick:
# Using Linear Matrix Inequality

$$\|A(\boldsymbol{\alpha})\|_{\infty} = \max \sigma \left[A(\boldsymbol{\alpha})\right] \leq \lambda$$

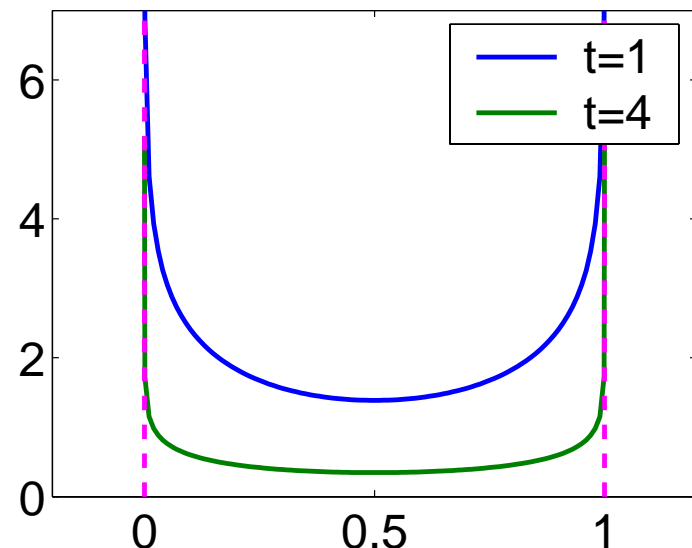$$A(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i y_i X_i$$

# The Second Trick:
# Using Linear Matrix Inequality

$$\|A(\boldsymbol{\alpha})\|_\infty = \max \sigma \left[ A(\boldsymbol{\alpha}) \right] \leq \lambda$$

$$A(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i y_i X_i$$

$$\Leftrightarrow \begin{bmatrix} \lambda I_R & A(\boldsymbol{\alpha}) \\ A^\top(\boldsymbol{\alpha}) & \lambda I_C \end{bmatrix} \succeq 0$$

Ryota Tomioka at ICML 2007 Corvallis, OR, USA

# The Third Trick:
# Interior Point Method

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad \sum_{i=1}^{n} \ell_{\mathsf{LR}}^*(\alpha_i) + \frac{1}{t}\phi(\boldsymbol{\alpha}),$$

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i y_i = 0.$$



$$t \rightarrow \infty$$

Original problem!

$$\phi(\boldsymbol{\alpha}) := -\left( \log \det \begin{bmatrix} \frac{\lambda}{n}I & A(\boldsymbol{\alpha}) \\ A^\top(\boldsymbol{\alpha}) & \frac{\lambda}{n}I \end{bmatrix} \right.$$
$$\left. + \log \alpha + \log(1-\alpha) \right).$$

$$(A(\boldsymbol{\alpha}) = \sum_i \alpha_i y_i X_i)$$

# Good News for IP optimization

- Obtaining the Primal Variable:

$$\widehat{\boldsymbol{\alpha}}_t \quad : \text{solution at barrier parameter } t$$

$$W_{\widehat{\alpha}_t} = U_{\widehat{\alpha}_t} \text{diag}\left( \frac{2\lambda_c^{(\widehat{\alpha}_t)}}{t\left(\lambda^2 - \lambda_c^{(\widehat{\alpha}_t)^2}\right)} \right) V_{\widehat{\alpha}_t}^\top$$

$$\left( U_{\widehat{\alpha}_t} \Lambda_{\widehat{\alpha}_t} V_{\widehat{\alpha}_t}^\top := A(\widehat{\boldsymbol{\alpha}}_t) = \sum_{i=1}^n \widehat{\alpha}_{t,i} y_i X_i \right)$$

- Quality guarantee:

$$\text{Duality gap}(W_{\widehat{\alpha}_t}, \widehat{\boldsymbol{\alpha}}_t) \leq \frac{R + C + 2n}{t}$$

Ryota Tomioka at ICML 2007 Corvallis, OR, USA

# Application: Motor-imagery EEG Classification

# Single-trial EEG Classification

The Covariance EEG signal            Class Label

$$X = SS^{\mathsf{T}} \implies y \in \{+1, -1\}$$

$$C \times C$$

$$S = $$

$$C \times T$$

# Single-trial EEG Classification

The Covariance EEG signal          Class Label

$$X = SS^\mathsf{T} \implies y \in \{+1, -1\}$$

$C \times C$
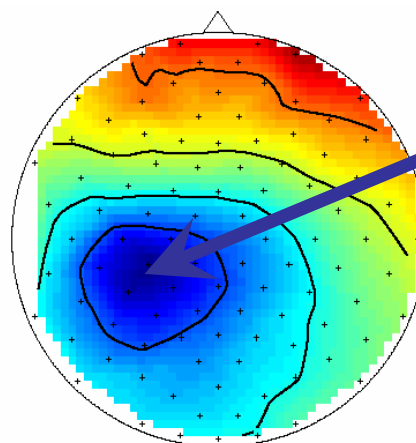
ERD/ERS

Lateralized modulation of rhythmic activity

C4 lap          Left                    Right          C3 lap

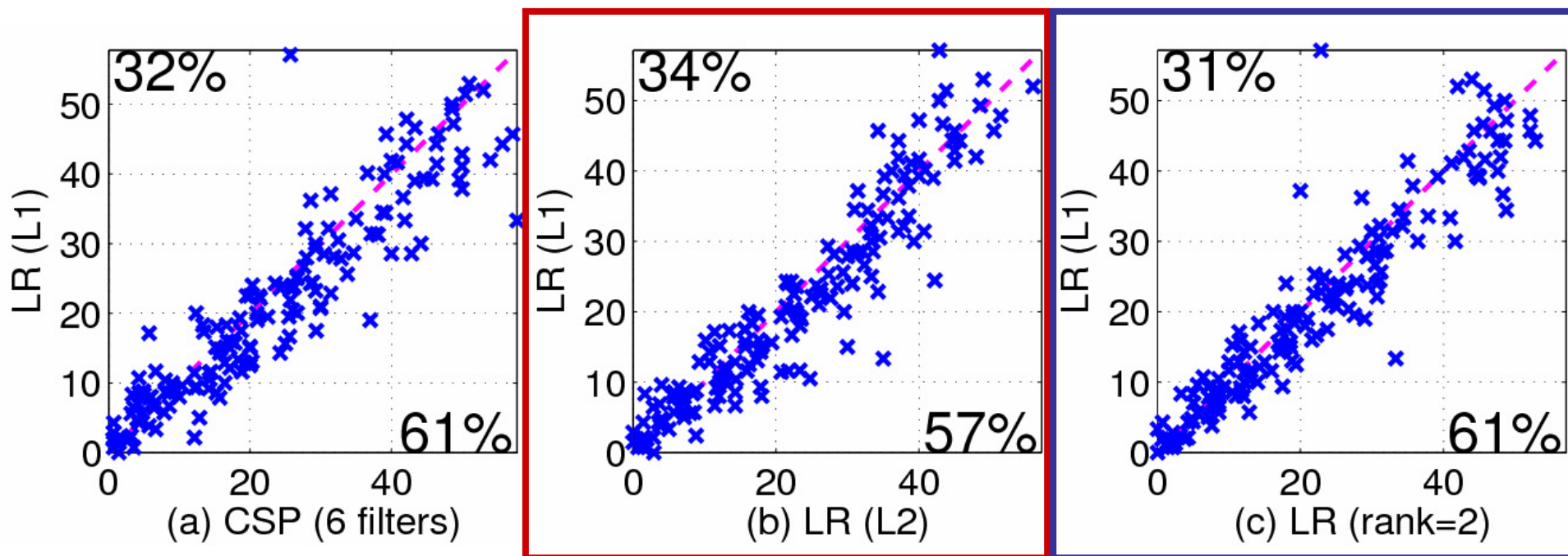# Conventional Methods

- Common Spatial Pattern (CSP) [Koles 1991; Ramoser 2000] (State of the art)
  - Two steps:
    - Feature Extraction: Find a low-dimensional decomposition.
    - Classify: linear classifier on the log-power feature.
- LR (L2)
  - $\ell_2$(Frobenius norm)-regularized logistic regression.
- LR (rank=2)
  - Rank=2 constrained logistic regression (nonconvex!)

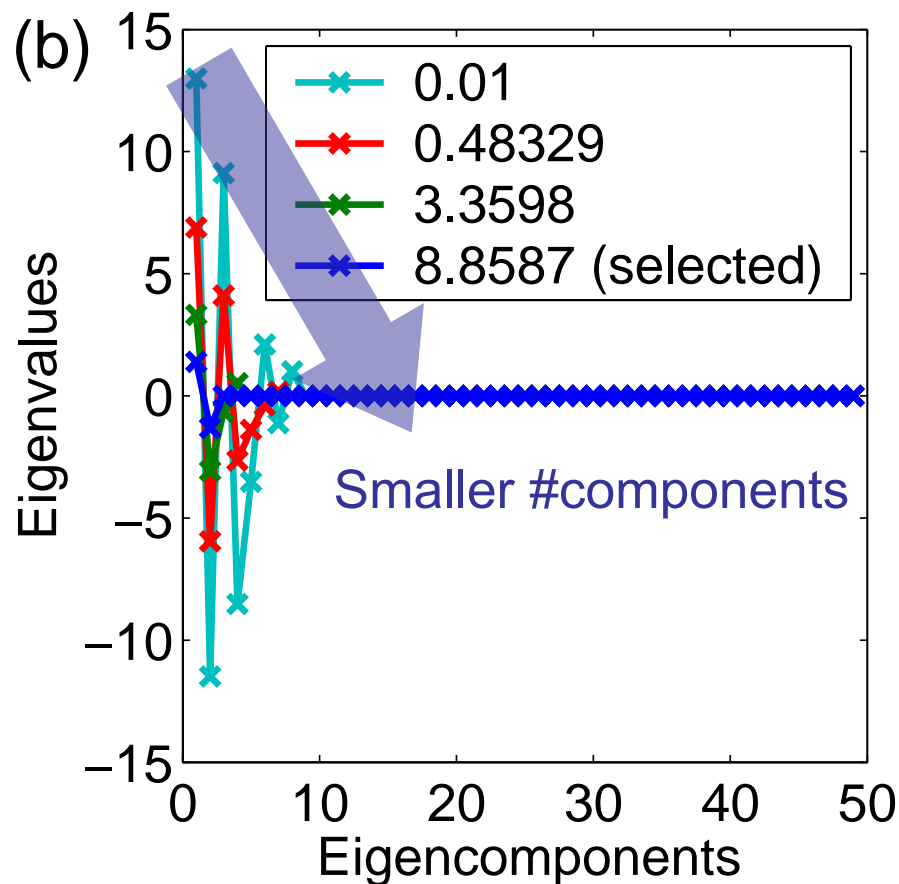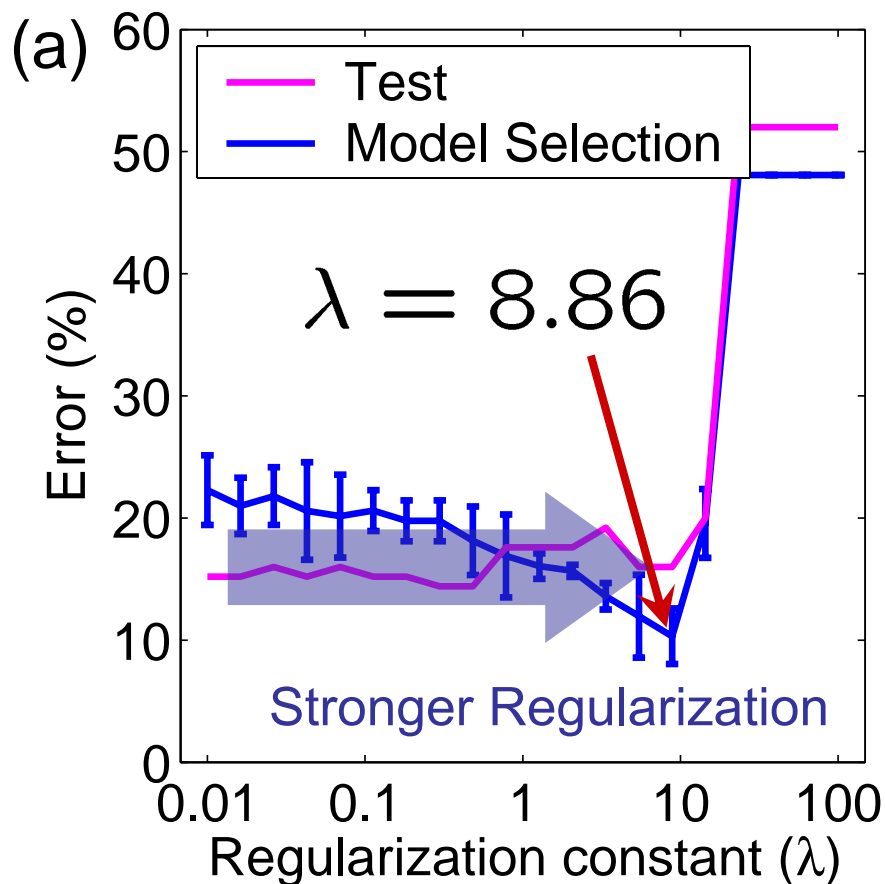$$W = \tfrac{1}{2}\left(-w_1 w_1^\top + w_2 w_2^\top\right)$$

# Results: Classification Errors



(a) CSP (6 filters)  (b) LR (L2)  (c) LR (rank=2)

- Low-ranked ($\ell_1$-regularized) solution performs better.
- Fixed rank performs suboptimal.

# Extracted Features (1/2)
# Model Selection and Eigenvalues



(a) Error (%) vs Regularization constant ($\lambda$). Legend: Test, Model Selection. $\lambda = 8.86$. Stronger Regularization.

(b) Eigenvalues vs Eigencomponents. Legend: 0.01, 0.48329, 3.3598, 8.8587 (selected). Smaller #components.

# Extracted Features (2/2)
# Eigenvectors



$$W = U \Lambda U^\top$$

$$\text{U}(:,1) \ (\lambda_1 = -1.31) \qquad \text{U}(:,2) \ (\lambda_1 = 1.40)$$

# Works on Spectral $\ell_1$ (Trace-norm) Regularization

- Prior work by Fazel, Hindi, and Boyd (2001)
- Related work by Abernethy et al. (2006)

|  | MMMF [Srebro et al. 05] | MTFL [Argyriou et al. 07] | Uncovering Shared Structure [Amit et al. 07] | Classifying Matrices [this talk] |
|---|---|---|---|---|
| Application | Matrix Factorization | Multi-ouput Regression | Multi-class Classification | Matrix Classification |
| Loss Function | Hinge-loss | Quad-loss | Hinge-loss | Logit-loss |
| Input | Scalar | Vector | Vector | Matrix |
| Output | Matrix | Vector | Vector | Scalar |
| Optimiza-tion | SDP | Iterative | Primal Gradient | Dual Interior-point |

# Summary

- Proposed the <span style="color:red">Matrix Classifier that factorizes</span> using the <span style="color:blue">Spectral $\ell_1$-regularization.</span>
  - Single <span style="color:blue">convex optimization</span> problem.
  - Dual formulation and Linear Matrix Inequality for efficient optimization.
  - <span style="color:blue">Sparseness</span>: interpretable solution.
- Applied to motor-imagery EEG classification
  - No distinction between feature extraction step and classification step.
  - Found physiologically relevant features.
  - Application to other problems are in progress.
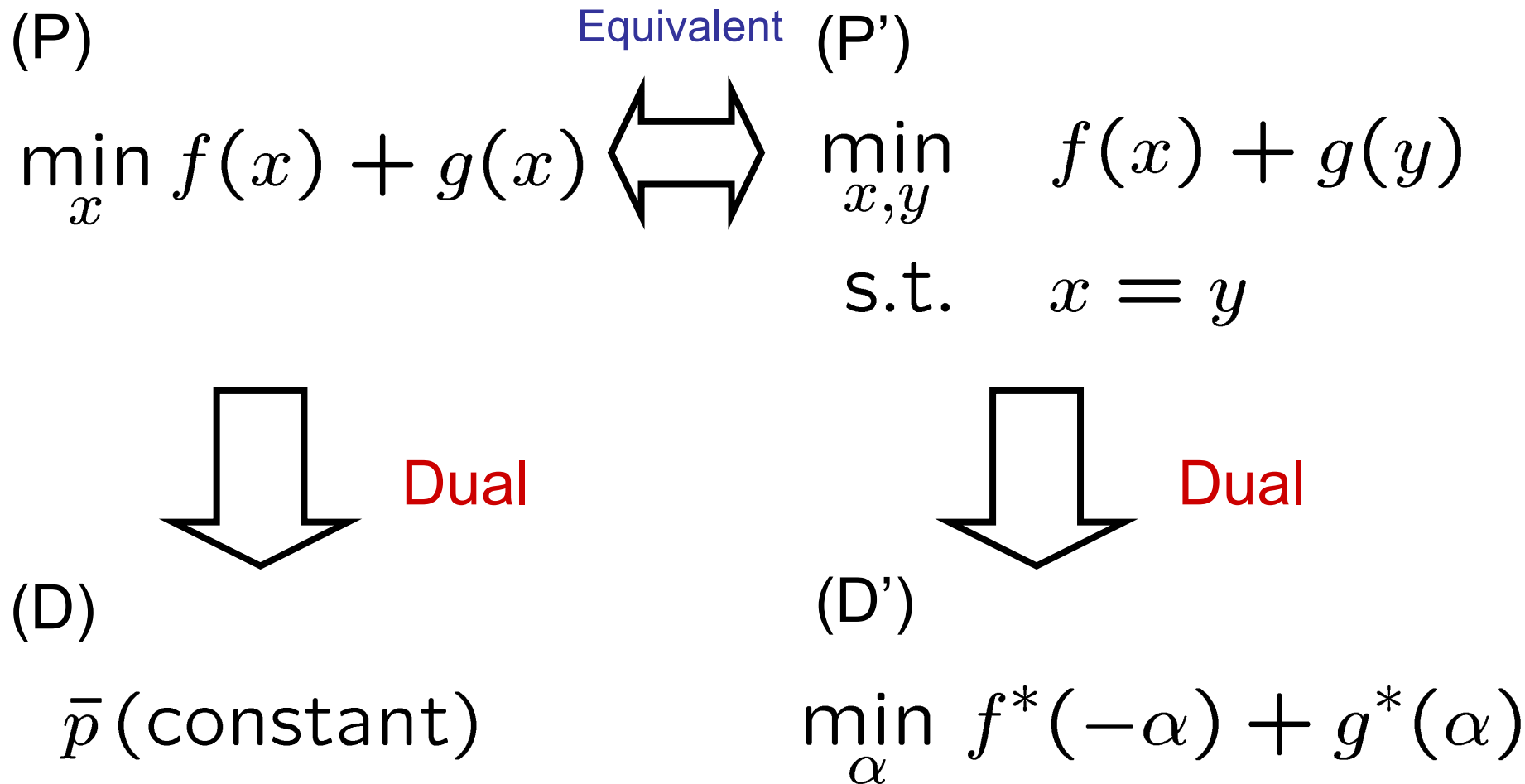
# References

- Koles (1991) "The quantitative extraction and topographic mapping of the abnormal components in the clinical EEG". *Electroencephalogr. Clin. Neurophysiol*., **79**.

- Ramoser et al. (2000) "Optimal spatial filtering of single trial EEG during imagined hand movement". *IEEE Trans. Rehab. Eng., **8**(4).

- Fazel et al. (2001). "A rank minimization heuristic with application to minimum order system approximation". *Proc. American Control Conference*.

- Boyd & Vandenberghe (2004). *Convex optimization*. CUP.

- Srebro et al. (2005) "Maximum margin matrix factorization". *Advances in NIPS.* **17**.

- Abernethy et al. (2006) , "Low-rank matrix factorization with attributes". *Technical report Ecole des Mines de Paris.* N24/06/MM.

- Blankertz et al. (2006) "The Berlin Brain-Computer Interface: EEG-based communication without subject training". *IEEE Trans. Neural Sys. Rehab. Eng*. **14**(2).

- Argyriou et al. (2007) "Multi-Task Feature Learning". *Advances in NIPS.* **19**.

- Tomioka et al. (2007) "Logistic regression for single trial EEG classification". *Advances in NIPS.* **19**.

- Amit et al. (2007) "Uncovering Shared Structures in Multiclass Classification". *Proc. ICML.*

# Thank you very much!

# Derivation of the Dual Problem

(P)    Equivalent    (P')

$$\min_{x} f(x) + g(x) \Longleftrightarrow \min_{x,y} \quad f(x) + g(y)$$

$$\text{s.t.} \quad x = y$$

Dual    Dual

(D)    (D')

$$\bar{p}\,(\text{constant}) \qquad \min_{\alpha} f^*(-\alpha) + g^*(\alpha)$$

# Derivation of the Dual

Logistic loss     $\ell_1$-norm

$$(P) \quad \min_{\substack{W \in \mathbb{R}^{R \times C}, b \in \mathbb{R}, \\ \boldsymbol{z} \in \mathbb{R}^n}} \quad \boxed{\frac{1}{n} \sum_{i=1}^{n} \ell_{LR}(z_i)} + \boxed{\frac{\lambda}{n} \|W\|_1},$$

$$\text{s.t.} \quad y_i \left( \text{Tr}\left[ W^\top X_i \right] + b \right) = z_i$$

$$(i = 1, \dots, n),$$

Dual logistic loss

$$(D) \quad \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad \boxed{\sum_{i=1}^{n} \ell^*_{\mathsf{LR}}(\alpha_i)\big|_{(0 \le \alpha_i \le 1)}} \quad \ell_\infty\text{-norm}$$

$$\sum_{i=1}^{n} \alpha_i y_i = 0, \quad \boxed{\left\| \sum_{i=1}^{n} \alpha_i y_i X_i \right\|_\infty} \le \lambda,$$

# Interpreting the dual variable

$$p_i = \begin{cases} 1 - \alpha_i & (y_i = +1) \\ \alpha_i & (y_i = -1) \end{cases} \quad (i = 1, \ldots, n)$$

(D) $\quad \max_{0 \leq \boldsymbol{p} \leq 1} \quad \sum_{i=1}^{n} H_2(p_i)$

s.t. $\quad \sum_{i=1}^{n} \left( y_i - \mathbb{E}[y_i | p_i] \right) = 0,$

$$\left\| \sum_{i=1}^{n} \left( y_i - \mathbb{E}[y_i | p_i] \right) X_i \right\|_{\infty} \leq 2\lambda,$$

# Experimental setup

- Offline analysis of 162 datasets from 29 healthy subjects recorded in the Berlin Brain Computer Interface (BBCI) project ([Blankertz et al., 2006], www.bbci.de).

- Binary classification of all the combinations of left hand (L), right hand (R), and foot (F) imaginary movement.

- Multi-channel EEG (32, 64, or 128ch) recordings (70-600 trials in a dataset).

- Band-pass filter 7-30Hz.

# Conventional Methods

- CSP (Koles, 1991; Ramoser, 2000) Dimensionality reduction/ demixing technique using label information:

$$\Sigma^{(+)}\boldsymbol{w}_c = \lambda_c\Sigma^{(-)}\boldsymbol{w}_c \quad (c = 1,\ldots,C)$$

$$\Sigma^{(\pm)} = \langle X \rangle_{\pm}$$

$$f(X) = \sum_{c=1}^{C^*} \beta_c \log\left[\boldsymbol{w}_c^{\top} X \boldsymbol{w}_c\right] + \beta_0$$

$$(C^* < C)$$

# Conventional Methods

- LR (L2) − Logistic regression with L2-regularization (Frobenius norm)

$$\Omega(W) = \tfrac{1}{2}\mathsf{Tr}\left[W^\top W\right]$$

- LR (rank=2) − Rank=2 approximated logistic regression

$$W = \tfrac{1}{2}\left(-\boldsymbol{w}_1\boldsymbol{w}_1^\top + \boldsymbol{w}_2\boldsymbol{w}_2^\top\right)$$

(Tomioka et al., NIPS*2006)

# The Discriminative Model

$$f(S; W, b) = \mathsf{Tr}\left[WSS^\top\right] + b$$



$$W \in \mathbb{S}^C$$

$$b \in \mathbb{R}$$

# Appendix: CSP (1/3)

- ## Common Spatial Pattern (CSP) [Koles, 1991]
  - – discriminative dimensionality reduction technique.



➡ Generalized eigenvalue problem
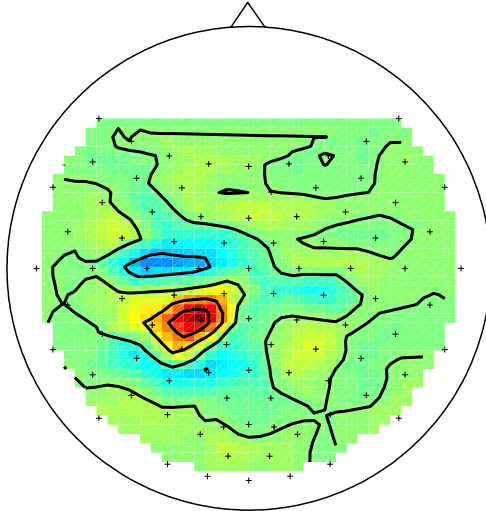
$$\Sigma^+ W = \Sigma^- W \Lambda$$
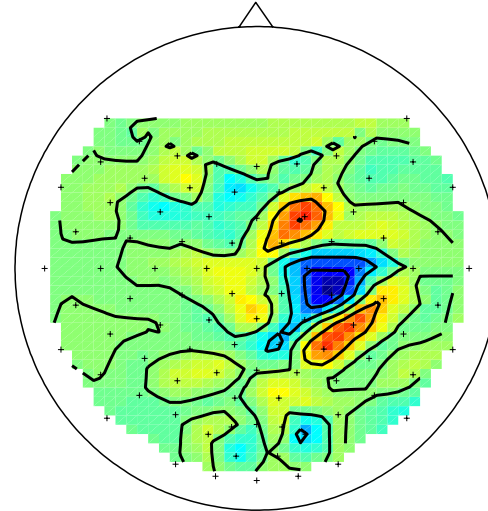
# Appendix: CSP (2/3)
## Example of Spatial Filters



left (−)    right (+)

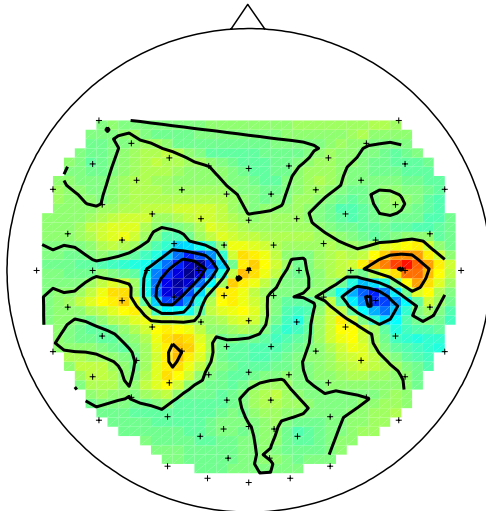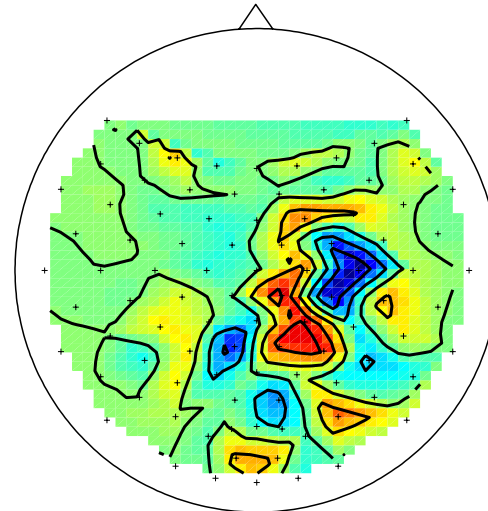csp1 [0.30]    csp3 [0.62]

csp2 [0.34]    csp4 [0.59]

# Appendix: CSP (3/3)
# CSP filtered time-series